

UNITED STATES PATENT APPLICATION

FOR

MULTI-THREADED SCHEDULED RECEIVE FOR FAST NETWORK PORT  
DATA

INVENTOR:

DONALD F. HOOPER  
MATTHEW J. ADILETTA  
GILBERT M. WOLRICH  
DEBRA BERNSTEIN

PREPARED BY:

CHARLES K. YOUNG  
INTEL CORPORATION  
2200 MISSION COLLEGE BLVD.  
SANTA CLARA, CA 95052-8119

(408) 765-8080

Attorney's Docket No. 42390.P7876X

"Express Mail" mailing label number EL23421818343

Date of Deposit 7/27/00

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to the Commissioner of Patents and Trademarks, Washington, D.C. 20231.

Kristin Baker  
(Typed or printed name of person mailing paper or fee)

Kristin Baker  
(Signature of person mailing paper or fee)

0966535-0700

# MULTI-THREADED SCHEDULED RECEIVE FOR FAST NETWORK PORT DATA

5

## BACKGROUND OF THE INVENTION

### 1. Field of the Invention

The described invention relates to the field of network communications. In particular, the invention relates to a method for scheduling multiple threads to process  
10 incoming network data.

### 2. Description of Related Art

Networking products such as routers require high speed components for packet  
15 data movement, i.e., collecting packet data from incoming network device ports and queuing the packet data for transfer to appropriate forwarding device ports. They also require high-speed special controllers for processing the packet data, that is, parsing the data and making forwarding decisions. Because the implementation of these high-speed functions usually involves the development of ASIC or custom devices, such  
20 networking products are of limited flexibility. For example, each controller is assigned to service network packets from one or more given ports on a permanent basis.

00626535-072700

## BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 is a block diagram of a communication system employing a hardware-based multi-threaded processor.

5           FIG. 2 is a block diagram of a microengine employed in the hardware-based multi-threaded processor of FIG. 1.

FIG. 3 is an illustration of an exemplary thread task assignment.

FIG. 4 is a block diagram of an I/O bus interface shown in FIG. 1.

10          FIG. 5 is a detailed diagram of a bus interface unit employed by the I/O bus interface of FIG. 4.

FIGS. 6A-6F are illustrations of various bus configuration control and status registers (CSRs).

FIG. 7A is a detailed diagram illustrating the interconnection between a plurality of 10/100 Ethernet (“slow”) ports and the bus interface unit.

15          FIG. 7B is a detailed diagram illustrating the interconnection between two Gigabit Ethernet (“fast”) ports and the bus interface unit.

FIGS. 8A-8C are illustrations of the formats of the RCV\_\_RDY\_CTL, RCV\_RDY\_HI and RCV\_RDY\_LO CSR registers, respectively.

20          FIG. 9 is a depiction of the receive threads and their interaction with the I/O bus interface during a receive process.

FIGS. 10A and 10B are illustrations of the format of the RCV\_REQ FIFO and the RCV\_CTL FIFO, respectively.

FIG. 11 is an illustration of the thread done registers.

25          FIG. 12 shows a graphical overview of the scheduling process showing an exemplary receive scheduler and three receive threads.

FIG. 13 shows a flowchart of the interaction of processes between the scheduler thread and an available receive thread .

FIG. 14 shows a flowchart showing the interaction of processing by multiple receive threads.

00/2/0" 5352960

## DETAILED DESCRIPTION

Referring to FIG. 1, a communication system 10 includes a parallel, hardware-based multi-threaded processor 12. The hardware based multi-threaded processor 12 is coupled to a first peripheral bus (shown as a PCI bus) 14, a second peripheral bus referred to as an I/O bus 16 and a memory system 18. The system 10 is especially useful for tasks that can be broken into parallel subtasks or functions. The hardware-based multi-threaded processor 12 includes multiple microengines 22, each with multiple hardware controlled program threads that can be simultaneously active and independently work on a task. In the embodiment shown, there are six microengines 22a-22f and each of the six microengines is capable of processing four program threads, as will be described more fully below.

The hardware-based multi-threaded processor 12 also includes a processor 23 that assists in loading microcode control for other resources of the hardware-based multi-threaded processor 12 and performs other general purpose computer type functions such as handling protocols, exceptions, extra support for packet processing where the microengines pass the packets off for more detailed processing. In one embodiment, the processor 23 is a StrongARM (ARM is a trademark of ARM Limited, United Kingdom) core based architecture. The processor (or core) 23 has an operating system through which the processor 23 can call functions to operate on the microengines 22a-22f. The processor 23 can use any supported operating system, preferably real-time operating system. For the core processor implemented as a StrongARM architecture, operating systems such as MicrosoftNT real-time, VXWorks and :CUS, a freeware operating system available over the Internet, can be used.

The six microengines 22a-22f each operate with shared resources including the memory system 18, a PCI bus interface 24 and an I/O bus interface 28. The PCI bus

interface provides an interface to the PCI bus 14. The I/O bus interface 28 is responsible for controlling and interfacing the processor 12 to the I/O bus 16. The memory system 18 includes a Synchronous Dynamic Random Access Memory (SDRAM) 18a, which is accessed via an SDRAM controller 26a, a Static Random Access Memory (SRAM) 18b, which is accessed using an SRAM controller 26b, and a nonvolatile memory (shown as a FlashROM) 18c that is used for boot operations. The SDRAM 16a and SDRAM controller 26a are typically used for processing large volumes of data, e.g., processing of payloads from network packets. The SRAM 18b and SRAM controller 26b are used in a networking implementation for low latency, fast access tasks, e.g., accessing look-up tables, memory for the processor 23, and so forth. The microengines 22a-22f can execute memory reference instructions to either the SDRAM controller 26a or the SRAM controller 18b.

The hardware-based multi-threaded processor 12 interfaces to network devices such as a media access controller device, including a “slow” device 30 (e.g., 10/100BaseT Ethernet MAC) and/or a “fast” device 31, such as Gigabit Ethernet MAC, ATM device or the like, over the I/O Bus 16. In the embodiment shown, the slow device 30 is an 10/100 BaseT Octal MAC device and thus includes 8 slow ports 32a-32h, and the fast device is a Dual Gigabit MAC device having two fast ports 33a, 33b. Each of the network devices attached to the I/O Bus 16 can include a plurality of ports to be serviced by the processor 12. Other devices, such as a host computer (not shown), that may be coupled to the PCI bus 14 are also serviced by the processor 12. In general, as a network processor, the processor 12 can interface to any type of communication device or interface that receives/sends large amounts of data. The processor 12 functioning as a network processor could receive units of packet data from the devices 30, 31 and process those units of packet data in a parallel manner, as will be described.

The unit of packet data could include an entire network packet (e.g., Ethernet packet) or a portion of such a packet.

Each of the functional units of the processor 12 are coupled to one or more internal buses. The internal buses include an internal core bus 34 (labeled “AMBA”) for coupling the processor 23 to the memory controllers 26a, 26b and to an AMBA translator 36. The processor 12 also includes a private bus 38 that couples the microengines 22a-22f to the SRAM controller 26b, AMBA translator 36 and the Fbus interface 28. A memory bus 40 couples the memory controllers 26a, 26b to the bus interfaces 24, 28 and the memory system 18.

Referring to FIG. 3, an exemplary one of the microengines 22a-22f is shown. The microengine 22a includes a control store 70 for storing a microprogram. The microprogram is loadable by the central processor 20. The microengine 70 also includes control logic 72. The control logic 72 includes an instruction decoder 73 and program counter units 72a-72d. The four program counters are maintained in hardware. The microengine 22a also includes context event switching logic 74. The context event switching logic 74 receives messages (e.g., SEQ\_#\_EVENT\_RESPONSE; FBI\_EVENT\_RESPONSE; SRAM\_EVENT\_RESPONSE; SDRAM\_EVENT\_RESPONSE; and AMBA\_EVENT\_RESPONSE) from each one of the share resources, e.g., SRAM 26b, SDRAM 26a, or processor core 20, control and status registers, and so forth. These messages provides information on whether a requested function has completed. Based on whether or not the function requested by a thread has completed and signaled completion, the thread needs to wait for that complete signal, and if the thread is enabled to operate, then the thread is placed on an available thread list (not shown). As earlier mentioned, in one embodiment, the microengine 22a can have a maximum of four threads of execution available.

In addition to event signals that are local to an executing thread, the microengine employs signaling states that are global. With signaling states, an executing thread can broadcast a signal state to all microengines 22. Any and all threads in the microengines can branch on these signaling states. These signaling states can be used to determine availability of a resource or whether a resource is due for servicing.

The context event logic 74 has arbitration for the four threads. In one embodiment, the arbitration is a round robin mechanism. However, other arbitration techniques, such as priority queuing or weighted fair queuing, could be used. The microengine 22a also includes an execution box (EBOX) data path 76 that includes an arithmetic logic unit (ALU) 76a and a general purpose register (GPR) set 76b. The ALU 76a performs arithmetic and logical functions as well as shift functions.

The microengine 22a further includes a write transfer register file 78 and a read transfer register file 80. The write transfer register file 78 stores data to be written to a resource. The read transfer register file 80 is for storing return data from a resource. Subsequent to or concurrent with the data arrival, an event signal from the respective shared resource, e.g., memory controllers 26a, 26b, or core 23, will be provided to the context event arbiter 74, which in turn alerts the thread that the data is available or has been sent. Both transfer register files 78, 80 are connected to the EBOX 76 through a data path. In the described implementation, each of the register files includes 64 registers.

The functionality of the microengine threads is determined by microcode loaded (via the core processor) for a particular user's application into each microengine's control store 70. Referring to FIG. 3, an exemplary thread task assignment 90 is shown. Typically, one of the microengine threads is assigned to serve as a receive scheduler 92 and another as a transmit scheduler 94. A plurality of threads



are configured as receive processing threads 96 and transmit processing (or “fill”) threads 98. Other thread task assignments include a transmit arbiter 100 and one or more core communication threads 102. Once launched, a thread performs its function independently.

5           The receive scheduler thread 92 assigns packets to receive processing threads 96. In a packet forwarding application for a bridge/router, for example, the receive processing thread parses packet headers and performs lookups based in the packet header information. Once the receive processing thread or threads 96 has processed the packet, it either sends the packet as an exception to be further processed  
10 by the core 23 (e.g., the forwarding information cannot be located in lookup and the core processor must learn it), or stores the packet in the SDRAM and queues the packet in a transmit queue by placing a packet link descriptor for it in a transmit queue associated with the transmit (forwarding port) indicated by the header/lookup. The transmit queue is stored in the SRAM. The transmit arbiter thread 100 prioritizes the  
15 transmit queues and the transmit scheduler thread 94 assigns packets to transmit processing threads that send the packet out onto the forwarding port indicated by the header/lookup information during the receive processing.

          The receive processing threads 96 may be dedicated to servicing particular ports or may be assigned to ports dynamically by the receive scheduler thread 92. For  
20 certain system configurations, a dedicated assignment may be desirable. For example, if the number of ports is equal to the number of receive processing threads 96, then it may be quite practical as well as efficient to assign the receive processing threads to ports in a one-to-one, dedicated assignment. In other system configurations, a dynamic assignment may provide a more efficient use of system resources.

25           The receive scheduler thread 92 maintains scheduling information 104 in the GPRs 76b of the microengine within which it executes. The scheduling information

104 includes thread capabilities information 106, port-to-thread assignments (list) 108 and “thread busy” tracking information 110. At minimum, the thread capabilities information informs the receive scheduler thread as to the type of tasks for which the other threads are configured, e.g., which threads serve as receive processing threads.

- 5 Additionally, it may inform the receive scheduler of other capabilities that may be appropriate to the servicing of a particular port. For instance, a receive processing thread may be configured to support a certain protocol, or a particular port or ports. A current list of the ports to which active receive processing threads have been assigned by the receive scheduler thread is maintained in the thread-to-port assignments list 108.
- 10 The thread busy mask register 110 indicates which threads are actively servicing a port. The receive scheduler uses all of this scheduling information in selecting threads to be assigned to ports that require service for available packet data, as will be described in further detail below.

- Referring to FIG. 4, the I/O bus interface 28 includes shared resources 120,
- 15 which are coupled to a push/pull engine interface 122 and a bus interface unit 124. The bus interface unit 124 includes a ready bus controller 126 connected to a ready bus 128 and an Fbus controller 130 for connecting to a portion of the I/O bus referred to as an Fbus 132. Collectively, the ready bus 128 and the Fbus 132 make up the signals of the I/O bus 16 (FIG. 1). The resources 120 include two FIFOs, a transmit FIFO 134 and a
- 20 receive FIFO 136, as well as CSRs 138, a scratchpad memory 140 and a hash unit 142. The Fbus 132 transfers data between the ports of the devices 30, 31 and the I/O bus interface 28. The ready bus 128 is an 8-bit bus that performs several functions. It is used to read control information about data availability from the devices 30, 31, e.g., in the form of ready status flags. It also provides flow control information to the devices
- 25 30, 31, and may be used to communicate with another network processor 12 that is connected to the Fbus 132. Both buses 128, 132 are accessed by the microengines 22



registers, including a thread status registers. Those of the registers which pertain to the receive process will be described in further detail.

The interrupt/signal registers include an INTER\_THD\_SIG register for inter-thread signaling. Any thread within the microengines 22 or the core 23 can write a thread number to this register to signal an inter-thread event.

Further details of the Fbus controller 130 and the ready bus controller 126 are shown in FIG. 5. The ready bus controller 126 includes a programmable sequencer 160 for retrieving MAC device status information from the MAC devices 30, 31, and asserting flow control to the MAC devices over the ready bus 128 via ready bus interface logic 161. The Fbus controller 130 includes Fbus interface logic 162, which is used to transfer data to and from the devices 30, 31, is controlled by a transmit state machine (TSM) 164 and a receive state machine (RSM) 166. In the embodiment herein, the Fbus 132 may be configured as a bidirectional 64-bit bus, or two dedicated 32-bit buses. In the unidirectional, 32-bit configuration, each of the state machines owns its own 32-bit bus. In the bidirectional configuration, the ownership of the bus is established through arbitration. Accordingly, the Fbus controller 130 further includes a bus arbiter 168 for selecting which state machine owns the Fbus 132.

Some of the relevant CSRs used to program and control the ready bus 128 and Fbus 132 for receive processes are shown in FIGS. 6A-6F. Referring to FIG. 6A, RDYBUS\_TEMPLATE\_PROGx registers 170 are used to store instructions for the ready bus sequencer. Each register of these 32-bit registers 170a, 170b, 170c, includes four, 8-bit instruction fields 172. Referring to FIG. 6B, a RCV\_RDY\_CTL register 174 specifies the behavior of the receive state machine 166. The format is as follows: a reserved field (bits 31:15) 174a; a fast port mode field (bits 14:13) 174b, which specifies the fast (Gigabit) port thread mode, as will be described; an auto push prevent window field (bits 12:10) 174c for specifying the autopush prevent window used by the

ready bus sequencer to prevent the receive scheduler from accessing its read transfer registers when an autopush operation (which pushes information to those registers) is about to begin; an autopush enable (bit 9) 174d, used to enable autopush of the receive ready flags; another reserved field (bit 8) 174e; an autopush destination field (bits 7:6) 174f for specifying an autopush operation's destination register; a signal thread enable field (bit 5) 174g which, when set, indicates the thread to be signaled after an autopush operation; and a receive scheduler thread ID (bits 4:0) 174h, which specifies the ID of the microengine thread that has been configured as a receive scheduler.

Referring to FIG. 6C, a REC\_FASTPORT\_CTL register 176 is relevant to receiving packet data from fast ports (fast port mode) only. It enables receive threads to view the current assignment of header and body thread assignments for the two fast ports, as will be described. It includes the following fields: a reserved field (bits 31:20) 176a; an FP2\_HDR\_THD\_ID field (bits 19:15) 176b, which specifies the fast port 2 header receive (processing) thread ID; an FP2\_BODY\_THD\_ID field (bits 14:10) 176c for specifying the fast port 2 body receive processing thread ID; an FP1\_HDR\_THD\_ID field (bits 9:5) 176d for specifying the fast port 1 header receive processing thread ID; and an FP1\_BODY\_THD\_ID field (bits 4:0) 176e for specifying the fast port 1 body processing thread ID. The manner in which these fields are used by the RSM 166 will be described in detail later.

Although not depicted in detail, other bus registers include the following: a RDYBUS\_TEMPLATE\_CTL register 178 (FIG. 6D), which maintains the control information for the ready bus and the Fbus controllers, for example, it enables the ready bus sequencer; a RDYBUS\_SYNCH\_COUNT\_DEFAULT register 180 (FIG. 6E), which specifies the program cycle rate of the ready bus sequencer; and an FP\_FASTPORT\_CTL register 182 (FIG. 6F), which specifies how many Fbus clock

cycles the RSM 166 must wait between the last data transfer and the next sampling of fast receive status, as will be described.

Referring to FIG. 7A, the MAC device 30 provides transmit status flags 200 and receive status flags 202 that indicate whether the amount of data in an associated transmit FIFO 204 or receive FIFO 206 has reached a certain threshold level. The ready bus sequencer 160 periodically polls the ready flags (after selecting either the receive ready flags 202 or the transmit ready flags 200 via a flag select 208) and places them into appropriate ones of the CSRs 138 by transferring the flag data over ready bus data lines 209. In this embodiment, the ready bus includes 8 data lines for transferring flag data from each port to the Fbus interface unit 124. The CSRs in which the flag data are written are defined as RCV\_RDY\_HI/LO registers 210 for receive ready flags and XMIT\_RDY\_HI/LO registers 212 for transmit ready flags, if the ready bus sequencer 160 is programmed to execute receive and transmit ready flag read instructions, respectively.

When the ready bus sequencer is programmed with an appropriate instruction directing it to interrogate MAC receive ready flags, it reads the receive ready flags from the MAC device or devices specified in the instruction and places the flags into RCV\_RDY\_HI register 210a and a RCV\_RDY\_LO register 210b, collectively, RCV\_RDY registers 210. Each bit in these registers corresponds to a different device port on the I/O bus.

Also, and as shown in FIG. 7B, the bus interface unit 124 also supports two fast port receive ready flag pins FAST\_RX1 214a and FAST\_RX2 214b for the two fast ports of the fast MAC device 31. These fast port receive ready flag pins are read by the RSM 166 directly and placed into an RCV\_RDY\_CNT register 216.

The RCV\_RDY\_CNT register 216 is one of several used by the receive scheduler to determine how to issue a receive request. It also indicates whether a flow control request is issued.

Referring to FIG. 8A, the format of the RCV\_RDY\_CNT register 216 is as follows: bits 31:28 are defined as a reserved field 216a; bit 27 is defined as a ready bus master field 216b and is used to indicate whether the ready bus 128 is configured as a master or slave; a field corresponding to bit 26 216c provides flow control information; bits 25 and 24 correspond to FRDY2 field 216d and FRDY1 field 216e, respectively. The FRDY2 216d and FRDY1 216e are used to store the values of the FAST\_RX2 pin 214b and FAST\_RX1 pin 214a, respectively, both of which are sampled by the RSM 166 each Fbus clock cycle; bits 23:16 correspond to a reserved field 216f; a receive request count field (bits 15:8) 216g specifies a receive request count, which is incremented after the RSM 166 completes a receive request and data is available in the RFIFO 136; a receive ready count field (bits 7:0) 216h specifies a receive ready count, an 8-bit counter that is incremented each time the ready bus sequencer 160 writes the ready bus registers RCV\_RDY\_CNT register 216, the RCV\_RDY\_LO register 210b and RCV\_RDY\_HI register 210a to the receive scheduler read transfer registers.

There are two techniques for reading the ready bus registers: “autopush” and polling. The autopush instruction may be executed by the ready bus sequencer 160 during a receive process (rxautopush) or a transmit process (txautopush). Polling requires that a microengine thread periodically issue read references to the I/O bus interface 28.

The rxautopush operation performs several functions. It increments the receive ready count in the RCV\_RDY\_CNT register 216. If enabled by the RCV\_RDY\_CTL register 174, it automatically writes the RCV\_RDY\_CNT 216, the RCV\_RDY\_LO and RCV\_RDY\_HI registers 210b, 210a to the receive scheduler read

transfer registers and signals to the receive scheduler thread 92 (via a context event signal) when the rxautopush operation is complete.

The ready bus sequencer 160 polls the MAC FIFO status flags periodically and asynchronously to other events occurring in the processor 12. Ideally, the rate at which the MAC FIFO ready flags are polled is greater than the maximum rate at which the data is arriving at the MAC ports. Thus, it is necessary for the receive scheduler thread 92 to determine whether the MAC FIFO ready flags read by the ready bus sequencer 160 are new, or whether they have been read already. The rxautopush instruction increments the receive ready count in the RCV\_RDY\_CNT register 216 each time the instruction executes. The RCV\_RDY\_CNT register 216 can be used by the receive scheduler thread 92 to determine whether the state of specific flags have to be evaluated or whether they can be ignored because receive requests have been issued and the port is currently being serviced. For example, if the FIFO threshold for a Gigabit Ethernet port is set so that the receive ready flags are asserted when 64 bytes of data are in the MAC receive FIFO 206, then the state of the flags does not change until the next 64 bytes arrive 5120 ns later. If the ready bus sequencer 160 is programmed to collect the flags four times each 5120 ns period, the next three sets of ready flags that are to be collected by the ready bus sequence 160 can be ignored.

When the receive ready count is used to monitor the freshness of the receive ready flags, there is a possibility that the receive ready flags will be ignored when they are providing new status. For a more accurate determination of ready flag freshness, the receive request count may be used. Each time a receive request is completed and the receive control information is pushed onto the RCV\_CNTL register 232, the the RSM 166 increments the receive request count. The count is recorded in the RCV\_RDY\_CNT register the first time the ready bus sequencer executes an rxrdy instruction for each program loop. The receive scheduler thread 92 can use this count to



track how many requests the receive state machine has completed. As the receive scheduler thread issues commands, it can maintain a list of the receive requests it submits and the ports associated with each such request.

Referring to FIGS. 8B and 8C, the registers RCV\_RDY\_HI 210a and  
5 RCV\_RDY\_LO 210b have a flag bit 217a, 217b, respectively, corresponding to each port.

Referring to FIG. 9, the receive scheduler thread 92 performs its tasks as quickly as possible to ensure that the RSM 166 is always busy, that is, that there is always a receive request waiting to be processed by the RSM 166. Several tasks  
10 performed by the receive scheduler 92 are as follows. The receive scheduler 92 determines which ports need to be serviced by reading the RCV\_RDY\_HI, RCV\_RDY\_LO and RCV\_RDY\_CNT registers 210a, 210b and 216, respectively. The receive scheduler 92 also determines which receive ready flags are new and which are old using either the receive request count or the receive ready count in the  
15 RCV\_RDY\_CNT register, as described above. It tracks the thread processing status of the other microengine threads by reading thread done status CSRs 240. The receive scheduler thread 92 initiates transfers across the Fbus 132 via the ready bus, while the receive state machine 166 performs the actual read transfer on the Fbus 132. The receive scheduler 92 interfaces to the receive state machine 166 through two FBI CSRs  
20 138: an RCV\_REQ register 230 and an RCV\_CNTL register 232. The RCV\_REQ register 230 instructs the receive state machine on how to receive data from the Fbus 132.

Still referring to FIG. 9, a process of initiating an Fbus receive transfer is shown. Having received ready status information from the RCV\_RDY\_HI/LO registers  
25 210a, 210b as well as thread availability from the thread done register 240 (transaction "1", as indicated by the arrow labeled 1), the receive scheduler thread 92 determines if

there is room in the RCV\_REQ FIFO 230 for another receive request. If it determines that RCV\_REQ FIFO 230 has room to receive a request, the receive scheduler thread 92 writes a receive request by pushing data into the RCV\_REQ FIFO 230 (transaction 2). The RSM 166 processes the request in the RCV\_REQ FIFO 230 (transaction 3). The

5 RSM 166 responds to the request by moving the requested data into the RFIFO 136 (transaction 4), writing associated control information to the RCV\_CTL FIFO 232 (transaction 5) and generating a start\_receive signal event to the receive processing thread 96 specified in the receive request (transaction 6). The RFIFO 136 includes 16 elements 241, each element for storing a 64 byte segment of data referred to herein as a

10 MAC packet ("MPKT"). The RSM 166 reads packets from the MAC ports in fragments equal in size to one or two RFIFO elements, that is, MPKTs. The specified receive processing thread 96 responds to the signal event by reading the control information from the RCV\_CTL register 232 (transaction 7). It uses the control information to determine, among other pieces of information, where the data is located

15 in the RFIFO 136. The receive processing thread 96 reads the data from the RFIFO 136 on quadword boundaries into its read transfer registers or moves the data directly into the SDRAM (transaction 8).

The RCV\_REQ register 230 is used to initiate a receive transfer on the Fbus and is mapped to a two-entry FIFO that is written by the microengines. The I/O bus

20 interface provides signals (not shown) to the receive scheduler thread indicating that the RCV\_REQ FIFO 236 has room available for another receive request and that the last issued receive request has been stored in the RCV\_REQ register 230.

Referring to FIG. 10A, the RCV\_REQ FIFO 230 includes two entries 231. The format of each entry 231 is as follows. The first two bits correspond to a reserved

25 field 230a. Bit 29 is an FA field 230b for specifying the maximum number of Fbus accesses to be performed for this request. A THSG field (bits 28:27) 230c is a two-bit



on the Fbus, where the data should be placed in the RFIFO and which microengine thread should be signaled once the data is received. The RSM 166 looks for a valid receive request in the RCV\_REQ FIFO 230. It selects the MAC device identified in the RM field and selects the specified port within the MAC by asserting the appropriate control signals. It then begins receiving data from the MAC device on the Fbus data lines. The receive state machine always attempts to read either eight or nine quadwords of data from the MAC device on the Fbus as specified in the receive request. If the MAC device asserts the EOP signal, the RSM 166 terminates the receive early (before eight or nine accesses are made). The RSM 166 calculates the total bytes received for each receive request and reports the value in the REC\_CNTL register 232. If EOP is received, the RSM 166 determines the number of valid bytes in the last received data cycle.

The RCV\_CNTL register 232 is mapped to a four-entry FIFO (referred to herein as RCV\_CNTL\_FIFO 232) that is written by the receive state machine and read by the microengine thread. The I/O bus interface 28 signals the assigned thread when a valid entry reaches the top of the RCV\_CNTL\_FIFO. When a microengine thread reads the RCV\_CNTL register, the data is popped off the FIFO. If the SIGRS field 230i is set in the RCV\_REQ register 230, the receive scheduler thread 92 specified in the RCV\_CNTL register 232 is signaled in addition to the thread specified in TID field 230k. In this case, the data in the RCV\_CNTL register 232 is read twice before the receive request data is retired from the RCV\_CTL\_FIFO 232 and the next thread is signaled. The receive state machine writes to the RCV\_CTL register 232 as long as the FIFO is not full. If the RCV\_CTL\_FIFO 232 is full, the receive state machine stalls and stops accepting any more receive requests.

Referring to FIG. 10B, the RCV\_CNTL\_FIFO 232 provides instruction to the signaled thread (i.e., the thread specified in TID) to process the data. As indicated

above, the RCV\_CNTL FIFO includes 4 entries 233. The format of the RCV\_CNTL FIFO entry 233 is as follows: a THMSG field (31:30) 23a includes the 2-bit message copied by the RSM from REC\_REQ register[28:27]. A MACPORT/THD field (bits 29:24) 232b specifies either the MAC port number or a receive thread ID, as will be described in further detail below. An SOP SEQ field (23:20) 232c is used for fast ports and indicates a packet sequence number as an SOP (start-of-packet) sequence number if the SOP was asserted during the receive data transfer and indicates an MPKT sequence number if SOP was not so asserted. An RF field 232d and RERR field 232e (bits 19 and 18, respectively) both convey receive error information. An SE field 232f (17:14) and an FE field 232g (13:10) are copies of the E2 and E1 fields, respectively, of the REC\_REQ. An EF field (bit 9) 232h specifies the number of RFIFO elements which were filled by the receive request. An SN field (bit 8) 232i is used for fast ports and indicates whether the sequence number specified in SOP\_SEQ field 232c is associated with fast port 1 or fast port 2. A VLD BYTES field (7:2) 232j specifies the number of valid bytes in the RFIFO element if the element contains in EOP MPKT. An EOP field (bit 1) 232k indicates that the MPKT is an EOP MPKT. An SOP field (bit 0) 232l indicates that the MPKT is an SOP MPKT.

FIG. 11 illustrates the format of the thread done registers 240 and their interaction with the receive scheduler and processing threads 92, 96, respectively, of the microengines 22. The thread done registers 240 include a first thread status register, TH\_DONE\_REG0 240a, which has 2-bit status fields 241a corresponding to each of threads 0 through 15. A second thread status register, TH\_DONE\_REG1 240b, has 2-bit status fields 241b corresponding to each of threads 16 through 23. These registers can be read and written to by the threads using a CSR instruction (or fast write instruction, described below). The receive scheduler thread can use these registers to determine which RFIFO elements are not in use. Since it is the receive scheduler thread

92 that assigns receive processing threads 96 to process the data in the RFIFO elements, and it also knows the thread processing status from the THREAD\_DONE\_REG0 and THREAD\_DONE\_REG1 registers 240a, 240b, it can determine which RFIFO elements are currently available.

5           The THREAD\_DONE CSRs 240 support a two-bit message for each microengine thread. The assigned receive thread may write a two-bit message to this register to indicate that it has completed its task. Each time a message is written to the THREAD\_DONE register, the current message is logically ORed with the new message. The bit values in the THREAD\_DONE registers are cleared by writing a "1",  
10 so the scheduler may clear the messages by writing the data read back to the THREAD\_DONE register. The definition of the 2-bit status field is determined in software. An example of four message types is illustrated in TABLE 1 below.

2-BIT MESSAGE	DEFINITION
00	Busy.
01	Idle, processing complete.
10	Not busy, but waiting to finish processing of entire packet.
11	Idle, processing complete for an EOP MPKT.

TABLE 1

The assigned receive processing threads write their status to the THREAD\_DONE register whenever the status changes. For example, a thread may immediately write 00 to the THREAD\_DONE register after the receive state machine signals the assigned thread. When the receive scheduler thread reads the THREAD\_DONE register, it can

look at the returned value to determine the status of each thread and then update its thread/port assignment list.

5 The microengine supports a fast\_wr instruction that improves performance when writing to a subset of CSR registers. The fast\_wr instruction does not use the push or pull engines. Rather, it uses logic that services the instruction as soon as the write request is issued to the FBI CSR. The instruction thus eliminates the need for the pull engine to read data from a microengine transfer register when it processes the command. The meaning of the 10-bit immediate data for some of the CSRs is shown below.

10

CSR	10-BIT IMMEDIATE DATA
INTER_THD_SIG	Thread number of the thread that is to be signaled.
THREAD_DONE	A 2-bit message that is shifted into a position relative to the thread that is writing the message.
THREAD_DONE_INCR1 THREAD_DONE_INCR2	Same as THREAD_DONE except that either the enqueue_seq1 or enqueue_seq2 is also incremented.
INCR_ENQ_NUM1 INCR_ENQ_NUM2	Write a one to increment the enqueue sequence number by one.

TABLE 2

15 It will be appreciated that the receive process as described herein assumes that no packet exemptions occurred, that is, that the threads are able to handle the packet processing without assistance from the core processor. Further, the receive process as described also assumes the availability of FIFO space. It will be appreciated that the various state machines must determine if there is room available in a FIFO, e.g., the RFIFO, prior to writing new entries to that FIFO. If a particular FIFO is full, the state

machine will wait until the appropriate number of entries has been retired from that FIFO.

#### Scheduling of Received Network Data

In one embodiment, one of the microengines 22a-f is designated as a receive scheduler 92 and three microengines 22a-f are designated to be receive threads 96.

Because there are four threads per microengine, there are twelve receive threads total. The receive scheduler 92 assigns a ready port to an available receive thread, and causes transfer of packet data into the network processor. The receive threads then analyze and store sections of incoming packet data (MPKTs). Since the receive threads can operate simultaneously, multiple portions of a MPKT can be processed by receive threads 96 at the same time.

FIG. 12 shows a graphical overview of the scheduling process showing an exemplary receive scheduler and three receive threads. The scheduling process will be described in more detail with respect to Figures 13 and 14. In Figure 12, the receive scheduler initiates three receive requests one after another (301). Each receive request initiates a bus transfer that wakes up an available thread (302). Thread X 303 is the first thread to wake up. It retrieves control information and saves state, as will be described in more detail later, and then continues to process its network data. Passing state information early on in the thread processing allows the network processor to maintain a high data rate since the threads can continue to process their packet section at their own slower pace.

Thread Y 304 wakes up from the second bus transfer. Similar to thread X, thread Y retrieves control information, saves state, then processes its network data. Thread Z 305 behaves similarly. If threads X, Y, and Z are in different microengines, then each thread may proceed without context switching relative to each other.



FIG. 13 shows a flowchart of the interaction of processes between the scheduler thread 92 and an available receive thread 96. At box 401, the receive scheduler 92 obtains port read status of the ports and thread available status to determine which ports and threads are available. The receive scheduler 92 selects one of the available threads (box 401). The receive scheduler 92 also selects a ready port (box 402). A port is ready if it has data in its receive buffer that meets some criteria (e.g., received number of bytes over a threshold, received a complete packet, etc). The receive scheduler 92 assembles the receive request, inserting the port number and ID of the receive thread to wake up (box 403).

The receive scheduler 92 increments a mailbox slot id (called MSEQ, for MPKT sequence) (box 404). The mailbox can be implemented in a register or fast memory. In one embodiment, the mailbox has sixteen slots and the receive scheduler 92 increments MSEQ then truncates it to four bits, ensuring that MSEQ will count from 0-15 then wrap around to 0 and repeat on successive iterations of the loop. The receive scheduler 92 inserts the MSEQ into an unused field of the receive request (box 405). The bus interface gets the request (box 406), then starts a hardware sequencer that selects the device/port and transfers data across the bus into a receive FIFO and wakes up the receive thread (box 410).

The receive thread, upon waking up (box 420), reads the receive control register and extracts the MSEQ value from it (box 421). The receive thread then restores information about the packet (box 422), using a test and clear operation on bit 31 of the first word in the MSEQ. Bit 31 set indicates valid information. Bit 31 is set to 1 by a previous thread, or by initialization. (Note that initially the first slot of MSEQ has bit 31 set to enable the very first thread to break out of the test and clear loop.) The receive thread then modifies the information as needed, and stores it to the mailbox slot MSEQ+1 (box 423), for the next thread to retrieve. The receive thread continues to

processes its network data (box 424), writes thread done and either swaps out context or waits on a new start receive (box 425).

FIG. 14 shows a flowchart showing the interaction of processing by multiple receive threads. The packet data is received on one or more ports that are connected to MAC devices, that are in turn connected to networks. The data is received as a series of MPKTs. For every receive request transaction the port bus interface transfers a section of the packet from a MAC device port to a receive FIFO in the network processor. The bus interface wakes up a designated receive thread. The receive thread then reads a control register (box 501) to obtain information about the bus transaction, i.e., which port, where in the receive FIFO the data is stored, size of data, type of port, message and other information. In the receive control register, there can be a cancel message that indicates that the port was not ready. If there is no cancel message, the receive thread restores the packet address from the mailbox slot indexed by MSEQ (box 502), increments it by the buffer size of the data (box 504), and saves it to the mailbox slot indexed by MSEQ+1 (box 505). After saving the updated address, the thread copies the data from RFIFO to DRAM (boxes 506, 520). The receive thread writes a thread done register to indicate it is available for re-assignment.

If there is a cancel (box 503), the receive thread does not increment the address, but merely restores it from mailbox indexed by MSEQ and saves it to mailbox indexed by MSEQ+1. In this case no data is copied to DRAM. The receive thread writes a thread done register to indicate it is available for re-assignment.

The sequence repeats, with the receive scheduler making requests and assigning mailbox slots, and receive threads processing the data. In one embodiment, the mailbox slot ids are assigned incrementally by the receive scheduler 92 and passed through the bus interface to the receive threads 96. They provide that successive receives from the same port will copy data to sequential buffers in DRAM.

The mailbox entry can be used to carry a variety of descriptive information used to process the packet. A few examples include: 1) it can carry a buffer freelist identifier that identifies which memory allocation freelist was used for the packet. This information is passed through the mailbox and ultimately to the transmit threads, which  
5 can return the buffer to the correct freelist; 2) it can carry an offset into the packet buffer, that will later be used by transmit threads to locate the beginning of the packet in the memory buffer. This is useful when packets are modified as part of a protocol conversion, resulting in a different packet size at transmit.

In one embodiment, the mailbox entry includes a relative buffer address that  
10 includes a sixteen bit value. The relative buffer address (or a scaled offset thereof) is added to a SRAM base address to provide a pointer to a packet link descriptor location in SRAM. Similarly, the relative buffer address (or a scaled offset thereof) is added to a SDRAM base address to provide a pointer to packet data location in SDRAM. The mailbox may also include such as information as an element count indicating the  
15 number of 64 byte segments of data in a data packet, and the SOP sequence number. The thread that processes an EOP MPKT writes the packet link descriptor to the right place. This may be to the transmit queue as previously described.

Thus, a method of processing network data in a network processor includes using a receive scheduler to schedule multiple threads. The scheduler allows greater  
20 flexibility by governing the rate at which each port or all ports receive data, causing selected discards. However, the specific arrangements and methods described herein are merely illustrative of the principles of this invention. Numerous modifications in form and detail may be made without departing from the scope of the described invention. Although this invention has been shown in relation to a particular  
25 embodiment, it should not be considered so limited. Rather, the described invention is limited only by the scope of the appended claims.